

# Detecting Unusual Input-Output Associations in Multivariate Conditional Data

Charmgil Hong\*

Milos Hauskrecht\*

## Abstract

Despite tremendous progress in outlier detection research in recent years, the majority of existing methods are designed only to detect *unconditional* outliers that correspond to unusual data patterns expressed in the joint space of all data attributes. Such methods are not applicable when we seek to detect *conditional* outliers that reflect unusual responses associated with a given context or condition. This work focuses on *multivariate conditional outlier detection*, a special type of the conditional outlier detection problem, where data instances consist of multi-dimensional input (context) and output (responses) pairs. We present a novel outlier detection framework that identifies abnormal input-output associations in data with the help of a decomposable conditional probabilistic model that is learned from all data instances. Since components of this model can vary in their quality, we combine them with the help of weights reflecting their reliability in assessment of outliers. We study two ways of calculating the component weights: global that relies on all data, and local that relies only on instances similar to the target instance. Experimental results on data from various domains demonstrate the ability of our framework to successfully identify multivariate conditional outliers.

## 1 Introduction

*Outlier detection* is a data analysis task that aims to find atypical behaviors, unusual outcomes, erroneous readings or annotations in data.<sup>1</sup> It has been an active research topic in data mining community, and it is frequently used in various applications to identify rare and interesting data patterns, which may be associated with beneficial or malicious events, such as fraud identification [27], network intrusion surveillance [10], disease outbreak detection [29], patient monitoring for preventable adverse events (PAE) [12, 11], *etc.* It is also utilized as a primary data preprocessing step that helps to remove noisy or irrelevant signals in data [14, 17].

Despite an extensive research, the majority of existing outlier methods are developed to detect *unconditional* outliers that are expressed in the joint space of all data attributes. Such methods may not work

well when one wants to identify *conditional* (contextual) outliers that reflect unusual responses for a given set of contextual attributes. Briefly, since conditional outliers depend on the context or properties of data instances, application of unconditional outlier detection methods may lead to incorrect results. For example, assume we want to identify incorrect (or highly unusual) image annotations in a collection of annotated images. Then by applying unconditional detection methods to the joint image-annotation space may lead to images with rare themes to be falsely identified as outliers due to the scarcity of these themes in the dataset, leading to false positives. Similarly, an unusual annotation of images with frequent themes may not be judged (scored) as very different from images with less frequent themes leading to false negatives.

This paper focuses on *multivariate conditional outlier detection*, a special type of the conditional outlier detection problem where data consists of  $m$ -dimensional continuous input vectors (context) and corresponding  $d$ -dimensional binary output vectors (responses). Our goal is to precisely identify the instances with unusual input-output associations. Following the definition of outlier given by Hawkins [13], we give a description of multivariate conditional outlier in plain language as:

**DEFINITION 1.** *A multivariate conditional outlier is an observation, which consists of context and associated responses, whose responses are deviating so much from the others in similar contexts as to arouse suspicions that it was generated by a different response mechanism.*

This formulation fits well various practical outlier detection problems that require contextual understanding of data. As briefly illustrated above, for example, recent social media services allow users to tag their content (*e.g.*, online documents, photos, or videos) with keywords and thereby permit keyword-based retrieval. These user annotations sometimes include irrelevant words by mistake that could be effectively pinpointed if the conditional relations between content and tags are considered. Likewise, evidence-based expert decisions (*e.g.*, functional categorization of genes, medical diagnosis and treatment decisions of patients) occasionally involve errors that could cause critical failures. Such erroneous decisions would be adequately detected through

\*Department of Computer Science, University of Pittsburgh.

<sup>1</sup>Outliers are also referred to as *anomalies*, *abnormalities*, *novelties*, *discordances*, or *deviants*.

contextual analysis of evidence-decision pairs.

The multivariate conditional outlier detection problem is challenging because both the contextual- and inter-dependences of data instances should be taken into account when identifying outliers. We tackle these challenges by building a probabilistic model  $P(\mathbf{Y}|\mathbf{X})$ , where  $\mathbf{X}=(X_1, \dots, X_m)$  denotes the input variables and  $\mathbf{Y}=(Y_1, \dots, Y_d)$  denotes the associated output variables. Briefly, the model is built (learned) from all available data, aiming to capture and summarize all relevant dependences among data attributes and their strength as observed in the data. Conditional outliers are then identified with the help of this model. More specifically, a conditional outlier corresponds to a data instance that is assigned a low probability by the model.

The exact implementation of the above approach is complicated, and multiple issues need to be resolved before it can be applied in practice. First, it is unclear how the probabilistic model  $P(\mathbf{Y}|\mathbf{X})$  should be represented and parameterized. To address this problem, we resort to and adapt structured probabilistic data models of  $P(\mathbf{Y}|\mathbf{X})$  that provide an efficient representation of input-output relations by decomposing the model using the chain rule into a product of univariate probabilistic factors  $P(Y_i|\mathbf{X}, \mathbf{Y}_{\pi(i)}) : i = 1, \dots, d$ ; i.e., each response  $Y_i$  is dependent on  $\mathbf{X}$  and a subset of the other responses  $\mathbf{Y}_{\pi(i)}$ . The univariate conditional models and their learning are rather common and well studied, and multiple models (e.g., logistic regression) can be applied to implement them. We note the structured probabilistic data models were originally proposed and successfully applied to support structured output prediction problems [30]. However, their application to outlier detection problems is new. The key difference is that while in prediction we seek to find outputs that maximize the probability given the inputs, in conditional outlier detection we aim to identify unusual (or low probability) associations in between observed inputs and outputs.

The second issue is that the probabilistic model must be learned from available data which can be hard especially when the number of context and output variables is high and the sample size is small. This may lead to model inaccuracies and miscalibration of probability estimates, which in turn may effect the identification of outliers. To alleviate this problem, we formulate and present outlier scoring methods that combine the probability estimates with the help of weights reflecting their reliability in assessment of outliers.

Through empirical studies, we test our approach on datasets with multi-dimensional responses. We demonstrate that our method is able to successfully identify multivariate conditional outliers and outperforms the existing baselines.

The rest of this paper is organized as follows. Section 2 formally define the problem. Section 3 reviews existing research on the topic. Section 4 describes our multivariate conditional outlier detection approach. Section 5 presents the experimental results and evaluations. Lastly, Section 6 summarizes the conclusions of our study.

## 2 Problem Definition

In this work, we study a special type of the conditional outlier detection problem where data consist of multi-dimensional input-output pairs; that is, each instance in dataset  $\mathcal{D} = \{\mathbf{x}^{(n)}, \mathbf{y}^{(n)}\}_{n=1}^N$  consists of an  $m$ -dimensional continuous input vector  $\mathbf{x}^{(n)} = (x_1^{(n)}, \dots, x_m^{(n)})$  and a  $d$ -dimensional binary output vector  $\mathbf{y}^{(n)} = (y_1^{(n)}, \dots, y_d^{(n)})$ . Our goal is to detect irregular response patterns in  $\mathbf{Y}$  given context  $\mathbf{X}$ . The fundamental issues in developing a multivariate conditional outlier detection method are how to take into account the *contextual dependences between output  $\mathbf{Y}$  and their input  $\mathbf{X}$* , as well as the *mutual dependences among  $\mathbf{Y}$* . We address these issues by building a decomposable probabilistic representation for  $\mathbf{Y}|\mathbf{X}$ .

Note that multivariate conditional outlier detection is clearly different from unconditional outlier detection when the problems are expressed probabilistically. In conditional outlier detection, we are interested in the instances that fall into low-probability regions of the conditional joint distribution  $P(\mathbf{y}|\mathbf{x}) = P(\mathbf{y}, \mathbf{x})/P(\mathbf{x})$ . On the other hand, unconditional outlier detection approaches generally seek instances in low-probability regions of the joint distribution  $P(\mathbf{y}, \mathbf{x})$ .

**Notation:** For notational convenience, we will omit the index superscript  $^{(n)}$  when it is not necessary. We may also abbreviate the expressions by omitting variable names; e.g.,  $P(Y_1 = y_1, \dots, Y_d = y_d|\mathbf{X} = \mathbf{x}) = P(y_1, \dots, y_d|\mathbf{x})$ .

## 3 Existing Research

Outlier detection has been extensively studied in the data mining and statistics communities [8, 15, 1]. A wide variety of approaches to tackle the detection problem for multivariate data have been proposed in the literature. Accordingly, depending on the type of outliers the method aims to detect, five general categories of *unconditional* outlier detection approaches appear in the literature. These include density-based approaches [5, 19], distance-based approaches [23, 4], depth-based approaches [21, 24], deviation-based approaches [3], and high-dimensional approaches [2, 15]. Below we briefly summarize each of these categories. For technical details, please refer to [8, 15, 1].

Density-based approaches assume that the density

around a normal data instance is similar to that of its neighbors [5, 19]. A typical representative method is Local Outlier Factor (LOF) [5], which measures a relative local density in  $k$ -nearest neighbor boundary. LOF has shown good performance in many applications and is considered as an off-the-shelf outlier detection method. In Section 5, we use LOF as the representative unconditional outlier detection method and compare the performance with our proposed approach.

Distance-based approaches assume that normal data instances come from dense neighborhoods, while outliers correspond to isolated points. A representative method is [23] which gives an outlier score to each instance using a robust variant of the Mahalanobis distance [22], measuring the distance between each instance to the main body of data distribution such that the instances located far from the center of data distribution are identified as outliers.

Depth-based approaches assume that outliers are at the fringe of the data regions and normal instances are close to or in the center of the region. The methods in this category assign depth  $k$  to each instance by gradually removing data from convex hulls, and the instances with small depth are considered as outliers [21]. A relevant method is the One-class Support Vector Machines [24], which assumes all the training data belong to the “normal” class and finds a decision boundary defining the region of normal data, whereas instances lie across the boundary are identified as outliers.

Deviation-based approaches assume that outliers are the outmost data instances in the data region and can be identified by measuring the impact of each instance on the variance of the dataset. One of the well-known algorithms in this category is Linear Method for Deviation Detection (LMDD) [3]. Compared to the depth-based approaches, deviation-based approaches do not require complicated contour generation process.

In high-dimensional spaces, the above approaches often fail because the distance metrics and density estimators become computationally intractable and analytically ineffective. Moreover, due to the sparsity of data, no meaningful neighborhood can be defined. High-dimensional approaches are proposed to handle such extreme cases. Typical methods in this category project the data to a lower dimensional subspace, such as grid-based subspace outlier detection [2]. For a detailed review on related methods, see [15].

While the vast majority of existing work were built to solve the unconditional outlier detection problem, the approaches may not work properly when it comes to *conditional* outliers, since they do not take into account the conditional relations among data attributes. Realizing this, recent years have seen increased interest in

the *conditional* outliers detection that aims to identify outliers in a set of outputs for given values of inputs. Several approaches have been proposed to address the problems in this regard [12, 11, 25]. However, these solutions either are limited to handle problems with a single output variable [12, 11] or assume a restricted relations among real-valued input and output variables through a Gaussian mixture [25]. As results, the existing methods either make an independence assumptions that is too restrictive or are unfit for modeling multi-dimensional binary output variables.

In contrast to the existing methods, our proposed approach is different in that (1) it properly models multi-label binary outputs by adopting a structured probabilistic data model to represent data; and (2) it utilizes the decomposed conditional probability estimates from individual response dimensions to identify outliers. Consequently, our proposed approach drives the process of outlier detection to a more granular level of the conditional behaviors in data and (as follows in Section 5) leads to a significant performance improvement in outlier detection. Furthermore, by maintaining separate models for individual output variables, our approach provides a practical advantage that the existing multivariate outlier detection methods do not allow. That is, one can delve into a trained multivariate conditional model and investigate the quality of each univariate representation  $P(y_i|\mathbf{x}, \mathbf{y}_{\pi(i)})$  to decide whether the individual model could be reliably used to support outlier detection. For example, a univariate model that produces inconsistent estimates could be preemptively excluded from the outlier detection phase. Since our goal is not to recover a complete data representation but to obtain a useful utility function for outlier detection, this sort of modularity allows us to utilize only the model with high confidence and, hence, to perform more robust outlier detection.

## 4 Our Approach

This section describes our approach to identify unusual input-output pairs, which we refer to as MCODE: *Multivariate Conditional Outlier DEtection*. To facilitate an effective detection method, we utilize a decomposable probabilistic data representation for  $P(\mathbf{Y}|\mathbf{X})$  to capture the dependence relations among inputs and outputs, and to assess outliers by seeking low-probability associations between them. Accordingly, having a precise probabilistic data model and proper outlier scoring methods is of primary concern. In Section 4.1, we discuss how to obtain an efficient data representation and accurate conditional probability estimates of observed input-output pairs, using the probabilistic structured data modeling approach [20]. In Section 4.2, we treat

the probability estimates as a proxy representation of observed instances and present two outlier scoring methods by analyzing the reliability of these estimates.

**4.1 Probabilistic Modeling and Estimation** Our MCODE approach works by analyzing data instances come in input-output pairs with a statistical model representing the conditional joint distribution  $P(\mathbf{Y}|\mathbf{X})$ . A direct learning of the conditional joint from data, however, is generally very expensive or even infeasible, because the number of possible output combinations grows exponentially with  $d$ . To avoid such a high cost of learning yet achieve an accurate data representation for outlier detection, we decompose the conditional joint into a product of conditional univariate distributions using the chain rule of probability:

$$(4.1) \quad P(Y_1, \dots, Y_d|\mathbf{X}) = \prod_{i=1}^d P(Y_i|\mathbf{X}, \mathbf{Y}_{\pi(i)})$$

where  $\mathbf{Y}_{\pi(i)}$  denotes the parents of  $Y_i$ ; *i.e.*, all the output variables preceding  $Y_i$  [20]. This decomposition lets us represent  $P(\mathbf{Y}|\mathbf{X})$  by simply specifying each univariate conditional factor,  $P(Y_i|\mathbf{X}, \mathbf{Y}_{\pi(i)})$ . In this work, we use a logistic regression model for each of the output dimensions, because it can effectively handle high-dimensional feature space defined by a mixture of continuous and discrete variables (*i.e.*,  $\mathbf{X}, \mathbf{Y}_{\pi(i)}$  conditioning  $Y_i$ ) using regularization [18, 7].<sup>2</sup>

In theory, the result of the above product should be invariant regardless of the chain order (order of  $Y_i$ ). Nevertheless, in practice, different chain orders produce different conditional joint distributions as they draw in models learned from different data [9]. For this reason, several structure learning methods that determine the optimal set of parents have been proposed [31, 16]. However, these methods require at least  $O(d^2 f_c)$  of time, where  $f_c$  denotes the time of learning a classifier, that would not be preferable, especially when the output dimensionality  $d$  is high.

In MCODE we address the above problem by relaxing the chain rule and by permitting circular dependencies among the output variables. That is, we let  $\mathbf{Y}_{\pi(i)}$ , the parents of  $Y_i$ , be all the remaining output variables, and assume the true dependence relations among them could be recovered through a proper regularization of logistic regression. To summarize, our structural decomposition allows us to capture the interactions among the output variables, as well as

the input-output relations, using a collection of individually trained probabilistic functions with a relaxed conditional independence assumption. We use  $\mathcal{M} = \{\theta_{\mathcal{M}(1)}, \dots, \theta_{\mathcal{M}(d)}\}$  to denote this structured data representation, where  $\theta_{\mathcal{M}(i)}$  is the parameters of the probabilistic model for the  $i$ -th output dimension. Assuming logistic regression, these base statistical functions are parameterized using  $\mathcal{D}$  as:

$$(4.2) \quad \theta_{\mathcal{M}(i)} = \arg \max_{\theta} \sum_{n=1}^N \log P(y_i^{(n)}|\mathbf{x}^{(n)}, \mathbf{y}_{-i}^{(n)}; \theta)$$

This defines a pseudo-conditional joint probability of an observation pair  $(\mathbf{x}, \mathbf{y})$  as:

$$(4.3) \quad \Psi(y_1, \dots, y_d|\mathbf{x}; \mathcal{M}) = \prod_{i=1}^d \tilde{P}(y_i|\mathbf{x}, \mathbf{y}_{-i}; \theta_{\mathcal{M}(i)})$$

where  $\mathbf{y}_{-i}$  denotes the values of all other output variables except  $Y_i$ .

Now let us apply our data representation  $\mathcal{M}$  to estimate the conditional probabilities of observed outputs. For notational convenience, we introduce an auxiliary vector  $\boldsymbol{\rho} = (\rho_1, \dots, \rho_d)$  of  $d$  random variables, each defined in a conditional probability space  $\rho_i = [0, 1]$ . Each element of  $\boldsymbol{\rho}$  is quantized by a probabilistic estimation process that is formalized as below by unleashing the product in Equation (4.3):

$$(4.4) \quad \mathcal{M} : (\mathbf{x}^{(n)}, \mathbf{y}^{(n)}) \rightarrow \boldsymbol{\rho}^{(n)} = (\rho_1^{(n)}, \dots, \rho_d^{(n)})$$

where

$$\rho_i^{(n)} = \begin{cases} \tilde{P}(y_i^{(n)}|\mathbf{x}^{(n)}, \mathbf{y}_{-i}^{(n)}; \theta_{\mathcal{M}(i)}) & \text{if } y_i^{(n)} = 1 \\ 1 - \tilde{P}(y_i^{(n)}|\mathbf{x}^{(n)}, \mathbf{y}_{-i}^{(n)}; \theta_{\mathcal{M}(i)}) & \text{otherwise.} \end{cases}$$

Accordingly, space of  $\boldsymbol{\rho}$  is projecting a normalized confidence level (*i.e.*, conditional probability estimate) of each observation  $(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$  across individual output dimensions, using the data representation  $\mathcal{M}$ . Figure 1 shows an illustrative example of this estimation where the input-output data instances (left) are projected to a 2-dimensional conditional probability space (right).

**4.2 Outlier Scoring** After the above probabilistic estimation process using  $\mathcal{M}$ , we consider the resultant conditional probabilities  $\boldsymbol{\rho}$  as proxies of the original instances, and further hypothesize that multivariate conditional outliers could be effectively detected in this proxy space where instances are analyzed and expressed in terms of univariate posterior probabilities. Our goal is now to define an *outlier score* that measures how unusual each input-output association is.

<sup>2</sup>Depending on data types and assumptions, different probabilistic classification functions could also be used, *e.g.*, naïve Bayes, relevance vector machine, or probabilistic support vector machine.

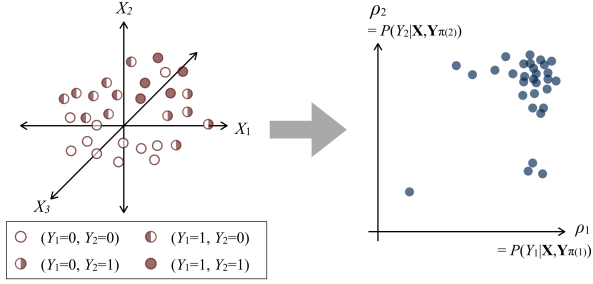


Figure 1: An illustrative example of conditional probability estimation process.

The most straightforward approach to define an outlier score is to use the probability  $P(\mathbf{y}|\mathbf{x})$  of data instances calculated by the model  $\mathcal{M}$ :

$$(4.5) \quad \text{Score}_{\text{PROD}}(\boldsymbol{\rho}^{(n)}) = - \sum_{i=1}^d \log \rho_i^{(n)}$$

Please note that this assumes all probability estimates and the models generating them are of high quality. However, in practice, the models that produce the probability estimates (*i.e.*,  $\theta_{\mathcal{M}(i)}$  in Equation (4.3)) may not be all equally reliable as they are trained from a finite number of samples (this is important especially when the number of input and output variables is high, and the sample size is small). Also, some dimensions of  $Y_i|\mathbf{X}, \mathbf{Y}_{\pi(i)}$  may not fit well the base statistical assumption (which in this work is a logistic curve) and result in miscalibrated estimations. Consequently, if we treat all dimensions of  $\boldsymbol{\rho}$  equally and merely search for the regions with low probabilities, the resulting scores degenerate to a noisy vector, which makes the detection of true irregularities hard.

To alleviate the issues, we propose to consider the reliability of each estimate dimension in  $\boldsymbol{\rho}$  (*i.e.*, the quality of model  $\theta_{\mathcal{M}(i)}$ ) and adjust their influence in outlier scoring by weights that reflect their reliability. We formalize our outlier score as:

$$(4.6) \quad \text{Score}_{\text{RW}}(\boldsymbol{\rho}^{(n)}) = - \sum_{i=1}^d w_i \log \rho_i^{(n)}$$

where  $w_i$  denotes the reliability weight of the model built for the  $i$ -th dimension. Note that, when  $w_i = 1$  for all dimensions  $i = 1, \dots, d$ , the score becomes equivalent to Equation (4.5), the negative log of the pseudo-conditional joint probability.

**4.2.1 Reliability Weights** One way to define reliability weights would be to use the Brier score [6] that measures the quality of the model in terms of model's

probability outputs. The Brier score is defined by averaging the squared errors of the probability estimates over all data instances:

$$\frac{1}{N} \sum_{n=1}^N (f^{(n)} - o^{(n)})^2$$

where  $f^{(n)}$  and  $o^{(n)}$  respectively denotes the predicted probability and actual outcome of the  $n$ -th instance. However, the assessment of the model quality for weighting purposes (Equation (4.6)) by the Brier score may not be the best as the score imposes different penalties for different errors (the mean squared error penalizes larger errors more than smaller errors) and varies the distribution of errors [28]. To address this, we propose our reliability weight be based on the mean estimated error, which gives the equal penalty to all errors:

**DEFINITION 2.** *Without loss of generality, let  $\epsilon_i^{(n)} = 1 - \rho_i^{(n)}$  be the estimated error of probability of an instance on dimension  $i$ . Reliability weight  $w_i$  is defined by taking the inverse of the mean estimated error:*

$$(4.7) \quad w_i = \frac{N}{\sum_{n=1}^N \epsilon_i^{(n)}}$$

Our variant of the Brier-like score estimates the quality of each estimate dimension  $\rho_i$  without distorting the distribution of errors. By taking the inverse of the score, we can effectively assign reliability weights to the dimensions, such that more on reliable dimensions become more important and the influence of noisy (unreliable) dimensions for outlier scoring is reduced.

**4.2.2 Local Reliability Weights** Notice that the above weighting scheme (Equation (4.7)) implicitly assumes that the reliability of probability estimates (*i.e.*, the quality of a model) is invariant across all data regions. However, the assumption often does not hold because in most practical problems especially with high-dimensional data spaces, data is not uniformly distributed in its attribute space. That is, modeling and estimation of  $P(Y_i|\mathbf{X}, \mathbf{Y}_{\pi(i)})$  cannot be achieved properly in sparse regions of the attribute space.

We tackle such a sparsity issue by evaluating the reliability of each dimension of  $\boldsymbol{\rho}$  locally in the region around the instance we want to check. This localized approach can be implemented as follows:

$$(4.8) \quad \text{Score}_{\text{LRW}}(\boldsymbol{\rho}^{(n)}) = - \sum_{i=1}^d w_i^{(n)} \log \rho_i^{(n)}$$

where

$$(4.9) \quad w_i^{(n)} = \frac{|N_k(n)|}{\sum_{n \in N_k(n)} \epsilon_i^{(n)}}$$

and  $N_k(n)$  denotes  $k$ -nearest neighbors of the  $n$ -th instance in the original attribute space. In the next section, we show the benefits of our reliability weights and outlier scores through experimental results.

## 5 Experiments

To validate and demonstrate the performance of our MCODE approach, we conduct experiments with data obtained from various domains. Through the empirical analysis in this section, we would like to verify the advantages of (1) adopting the conditional outlier detection approach, (2) considering the dependence relations among outputs, (3) weighting via reliability estimation, and (4) local reliability estimates and local outlier scores. Below we describe our experimental design and present the evaluation results.

**5.1 Compared Methods** To achieve our objectives, we perform experiments with the following methods:

- *Local outlier factor* (LOF) [5] – LOF is an unconditional method that estimates outliers using a relative local density measure in the joint space of all data attributes:

$$LOF((\mathbf{x}, \mathbf{y}), k) = \frac{\sum_{(\mathbf{x}', \mathbf{y}') \in N_k(\mathbf{x}, \mathbf{y})} \frac{lrd_k(\mathbf{x}', \mathbf{y}')}{lrd_k(\mathbf{x}, \mathbf{y})}}{|N_k(\mathbf{x}, \mathbf{y})|}$$

where  $N_k(\mathbf{x}, \mathbf{y})$  denotes the  $k$ -nearest neighborhood of instance  $(\mathbf{x}, \mathbf{y})$  and

$$lrd_k(\xi) = \frac{|N_k(\xi)|}{\sum_{o \in N_k(\xi)} \max(k\text{-dist}(o), \text{dist}(\xi, o))}$$

is the local reachability density which measures the geometric dispersion of the  $k$ -nearest neighborhood. LOF effectively finds the instances fall in sparse regions of data.

- *Conditional outlier detection with  $d$  independent response models* (I-PROD) – We apply [11] to the multivariate conditional setting by learning  $d$  independent conditional probability models  $P(Y_i|\mathbf{X})$  ( $Y_i$  is not dependent on other output variables) and scoring based on the product of their estimates (Equation (4.5)). We refer to this method as I-PROD.
- *MCODE without weighting* (M-PROD) (Equation (4.5))
- *MCODE with Reliability Weights* (M-RW) (Equation (4.6))

Dataset	$N/m/d$	Domain	Value Context	Description Response
Mediamill	43,907 / 120 / 101	Video	Video frames	Concepts
Enron	1,702 / 1,001 / 53	Text	Emails	Properties
Bibtex	7,395 / 1,836 / 159	Text	Paper metadata	Topics
Yahoo-business	11,214 / 21,924 / 30	Text	News articles	Topics
Yahoo-arts	7,484 / 23,146 / 26	Text	News articles	Topics
Yeast	2,417 / 103 / 14	Biology	Genes	Functionalities
Genbase	662 / 1,185 / 27	Biology	Genes	Functionalities
Birds	645 / 276 / 19	Sound	Bird songs	Species

Table 1: Dataset characteristics. ( $N$ : number of instances,  $m$ : input dimensionality,  $d$ : output dimensionality)

- *MCODE with Local Reliability Weights* (M-LRW) (Equation (4.8))

To obtain data models in I-PROD, M-PROD, M-RW, and M-LRW, we use  $L_2$ -penalized logistic regression and choose their regularization parameters by cross validation. In LOF and M-LRW, we set the number of neighbors  $k = 100$ .

**5.2 Data** We use *eight* public datasets with multi-dimensional input and output.<sup>3</sup> These are collected from various application domains, including sound recognition (*Birds*), biology (*Yeast*, *Genbase*), text categorization (*Yahoo* datasets, *Bibtex*, *Enron*), and semantic video/image annotation (*Mediamill*). Table 1 summarizes the characteristics of the datasets, such as dataset size, data domain, and short descriptions of the input and output variables.

**5.2.1 Simulating Outliers** For the purpose of our comparative evaluation, we simulate multivariate conditional outliers by perturbing the output space of data. There are two parameters in our simulation process: *Outlier ratio* specifies how many outliers per simulation are injected. We set this parameter to 1% throughout the experimental study. *Outlier dimensionality* specifies how many output dimensions of an outlier to be perturbed. We vary this parameter relative to the dimensionality of the output by perturbing {2.5, 5, 10, 20}% of outputs. To summarize, we simulate outliers as:

1. In each dataset, select 1% of instances uniformly at random
2. For each of the selected instances, perturb the values of {2.5, 5, 10, 20}% of the output dimensions (*i.e.*,  $y_{\text{perturbed}} = |y_{\text{original}} - 1|$ ) uniformly at random

We would like to stress that all methods (including their model building and detection stages) are always run on data with injected outliers. That is, we never learn a

<sup>3</sup>Datasets are available at <http://mulan.sourceforge.net> [26].

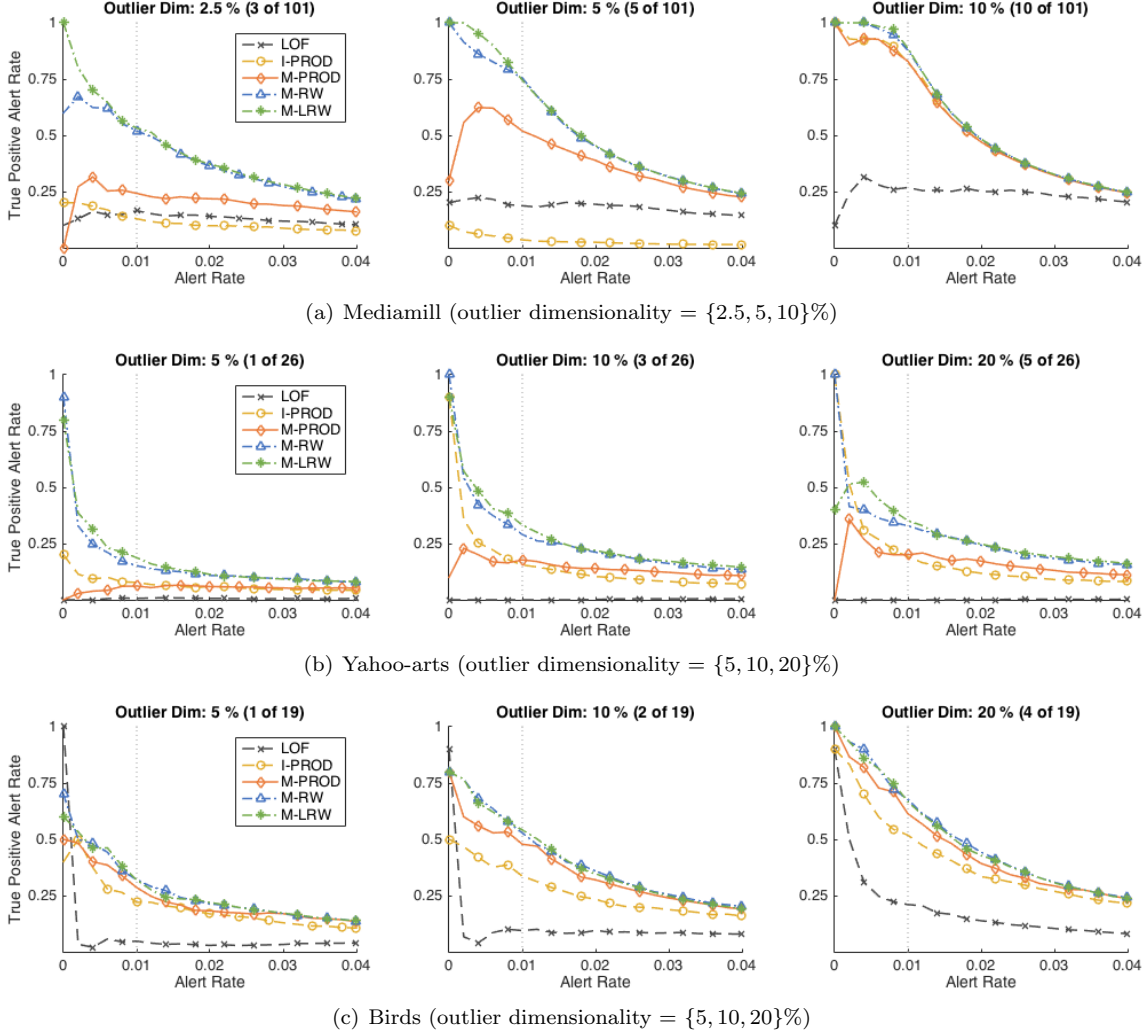


Figure 2: True positive alert rates at different alert rate ranging between 0 and 0.04.

model on the original (unperturbed) data and detect outliers on the simulated (perturbed) data. Such an design would be unrealistic since we do not know ahead of time what data instances to remove to learn a model.

Note that the simulated outliers can be analogous to the errors or mistakes in each application domain. For example, in semantic video/image annotation, perturbed output values can be perceived as inaccurate subject labels.

**5.3 Evaluation Metrics** We use *true positive alert rate* (TPAR) as our evaluation metric:

$$TPAR = (\text{True positive outliers}) / (\text{Predicted outliers})$$

TPAR (or precision) measures the percentage of instances with perturbation in the total number of instances detected by the methods. We assess TPAR in

two ways: We first evaluate TPAR at different alert rate (detection threshold) and analyze the quality of outlier scores (see Figure 2). We also measure the *Averaged TPAR* (ATPAR) in  $[0, 0.01]$  range, which coincides with the outlier ratio in our experiment setting. For both TPAR and ATPAR, higher is better.

**5.4 Results** Figure 2 and Table 2 show the performance of the five compared methods. All results are obtained from *ten* repeats.

Figures 2(a), 2(b), and 2(c) present the results on three datasets (*Mediamill*, *Yahoo-arts*, and *Birds*) for different outlier dimensions. Each figure illustrates the TPARs of all methods; X-axes show the alert rate, ranging between 0 and 0.04; Y-axes show TPAR. The vertical gray line at alert rate = 0.01 indicates where the alert rate is equal to the injected outlier ratio.

ATPAR	Outlier dimensionality = 2.5%						Outlier dimensionality = 5.0%					
	Baselines		MCODE				Baselines		MCODE			
	LOF	I-PROD	M-PROD	M-RW	M-LRW		LOF	I-PROD	M-PROD	M-RW	M-LRW	
Mediamill	0.14 ± 0.16	0.17 ± 0.09	0.26 ± 0.17	<b>0.61 ± 0.12</b>	<b>0.69 ± 0.09</b>		0.20 ± 0.17	0.06 ± 0.05	0.57 ± 0.14	<b>0.85 ± 0.05</b>	<b>0.90 ± 0.04</b>	
Enron	0.01 ± 0.03	0.12 ± 0.19	0.11 ± 0.11	0.06 ± 0.11	0.05 ± 0.09		0.01 ± 0.03	0.15 ± 0.22	0.20 ± 0.17	0.17 ± 0.22	0.21 ± 0.26	
Bibtex	0.00 ± 0.00	0.25 ± 0.28	0.32 ± 0.30	0.27 ± 0.29	0.33 ± 0.30		0.00 ± 0.01	<b>0.44 ± 0.27</b>	<b>0.49 ± 0.28</b>	<b>0.47 ± 0.28</b>	<b>0.51 ± 0.27</b>	
Yahoo-business	0.01 ± 0.02	0.13 ± 0.06	<b>0.21 ± 0.10</b>	<b>0.36 ± 0.09</b>	<b>0.38 ± 0.07</b>		0.01 ± 0.03	0.25 ± 0.08	<b>0.43 ± 0.11</b>	<b>0.56 ± 0.08</b>	<b>0.58 ± 0.07</b>	
Yahoo-arts	-	-	-	-	-		0.00 ± 0.01	0.11 ± 0.07	0.04 ± 0.04	<b>0.26 ± 0.06</b>	<b>0.29 ± 0.08</b>	
Genbase	-	-	-	-	-		0.05 ± 0.08	<b>0.93 ± 0.05</b>	<b>0.93 ± 0.06</b>	<b>0.94 ± 0.06</b>	<b>0.95 ± 0.06</b>	
Birds	-	-	-	-	-		0.04 ± 0.08	<b>0.34 ± 0.22</b>	<b>0.39 ± 0.25</b>	<b>0.45 ± 0.21</b>	<b>0.46 ± 0.22</b>	

ATPAR	Outlier dimensionality = 10.0%						Outlier dimensionality = 20.0%					
	Baselines		MCODE				Baselines		MCODE			
	LOF	I-PROD	M-PROD	M-RW	M-LRW		LOF	I-PROD	M-PROD	M-RW	M-LRW	
Mediamill	0.27 ± 0.16	<b>0.92 ± 0.03</b>	<b>0.91 ± 0.04</b>	<b>0.97 ± 0.03</b>	<b>0.98 ± 0.03</b>		0.30 ± 0.12	<b>0.99 ± 0.02</b>	<b>0.99 ± 0.01</b>	<b>1.00 ± 0.01</b>	<b>1.00 ± 0.00</b>	
Enron	0.03 ± 0.05	<b>0.21 ± 0.24</b>	<b>0.23 ± 0.17</b>	<b>0.31 ± 0.27</b>	<b>0.39 ± 0.29</b>		0.02 ± 0.04	0.26 ± 0.32	0.35 ± 0.21	<b>0.62 ± 0.24</b>	<b>0.78 ± 0.17</b>	
Bibtex	0.00 ± 0.01	<b>0.70 ± 0.20</b>	<b>0.70 ± 0.18</b>	<b>0.71 ± 0.20</b>	<b>0.72 ± 0.17</b>		0.00 ± 0.00	<b>0.88 ± 0.11</b>	<b>0.86 ± 0.12</b>	<b>0.88 ± 0.11</b>	<b>0.87 ± 0.11</b>	
Yahoo-business	0.01 ± 0.01	0.32 ± 0.10	<b>0.42 ± 0.13</b>	<b>0.57 ± 0.06</b>	<b>0.57 ± 0.07</b>		0.01 ± 0.02	<b>0.36 ± 0.13</b>	0.25 ± 0.09	<b>0.39 ± 0.05</b>	<b>0.41 ± 0.04</b>	
Yahoo-arts	0.00 ± 0.00	0.29 ± 0.05	0.18 ± 0.07	<b>0.44 ± 0.06</b>	<b>0.47 ± 0.05</b>		0.00 ± 0.00	0.36 ± 0.05	0.26 ± 0.07	0.39 ± 0.06	<b>0.45 ± 0.07</b>	
Yeast	0.08 ± 0.07	0.04 ± 0.06	0.45 ± 0.11	<b>0.64 ± 0.06</b>	<b>0.63 ± 0.05</b>		0.13 ± 0.09	0.17 ± 0.11	<b>0.52 ± 0.08</b>	<b>0.56 ± 0.07</b>	<b>0.54 ± 0.08</b>	
Genbase	0.06 ± 0.11	<b>0.96 ± 0.06</b>	<b>0.96 ± 0.04</b>	<b>0.98 ± 0.02</b>	<b>0.98 ± 0.02</b>		0.03 ± 0.09	<b>0.98 ± 0.03</b>	<b>0.96 ± 0.03</b>	<b>0.98 ± 0.02</b>	<b>0.98 ± 0.03</b>	
Birds	0.07 ± 0.11	<b>0.42 ± 0.31</b>	<b>0.56 ± 0.24</b>	<b>0.66 ± 0.18</b>	<b>0.66 ± 0.19</b>		0.32 ± 0.22	<b>0.67 ± 0.25</b>	<b>0.78 ± 0.19</b>	<b>0.85 ± 0.12</b>	<b>0.84 ± 0.13</b>	

Table 2: Averaged true positive alert rate in  $[0, 0.01]$ . Numbers shown in bold indicate the best results on each experiment set (by paired t-test at  $\alpha=0.05$ ).

In general, TPARs improve as the outlier dimensionality increases, because outliers with larger perturbations are easier to detect. Comparing the conditional outlier detection approaches (I-PROD, M-PROD, M-RW, and M-LRW) with the unconditional approach (LOF), the conditional approaches are clear winners as the conditional methods outperform LOF in most cases. This shows the advantages of the conditional outlier detection approaches in addressing the problem. Only exceptions are I-PROD on *Mediamill* when outlier dimensionality is low. This is because I-PROD does not consider the dependence relations among the output variables. Such advantages in modeling the inter-dependences of the outputs are consistently observed as M-PROD outperforms I-PROD in most experiments.

To show the benefits of our reliability weights, we analyze the performance of M-RW and M-LRW in comparison to that of M-PROD. An interesting point is that M-RW and M-LRW not only improve the performance drastically, but also make TPARs stable. This confirms that our reliability weighting methods can effectively estimate the quality of the models, and the resulting weights are useful in outlier scoring. Lastly, although M-LRW does not show much improvement from M-RW compared to the other key components of MCODE that we have discussed, the local weights seem to make M-RW even more stable as shown with *Mediamill* and *Yahoo-arts*.

Table 2 summarizes the results on all eight datasets in terms of ATPAR at 0.01. The table consists of four sections grouped by different values of outlier dimensionality ( $\{2.5, 5, 10, 20\}\%$ ). We do not report the results on the first four datasets (*Birds*, *Yeast*, *Genbase*, and *Yahoo-arts*) for outlier dimensionality

= 2.5% (for *Yeast*, 2.5% and 5.0%) because the output dimensionality ( $d$ ) is too small. The best performing methods on each experiment are shown in bold.

The results confirms the conclusions that we have drawn with Figure 2. One interesting point is that LOF shows exceptionally high (compared with its performance on other datasets) ATPAR on *Mediamill*. This is because the dataset has a similar number of input and output variables; hence, as outlier dimensionality increases, the simulated outliers become like unconditional outliers.

## 6 Conclusions

In this work, we introduced and tackled multivariate conditional outlier detection, a special type of the conditional outlier detection problem. We briefly reviewed existing research and motivated this new type of outlier detection problem. We presented our novel outlier detection framework that analyzes and detects abnormal input-output associations in data using a decomposable conditional probabilistic model that is learned from all data instances. We discussed how to obtain an efficient data representation and accurate conditional probability estimates of observed input-output pairs, using the probabilistic structured data modeling approach. Motivated by the Brier score, we developed present two outlier scoring methods by analyzing the reliability of probability estimates. Through the experimental results, we demonstrated the ability of our framework to successfully identify multivariate conditional outliers.

## References



- [1] Charu C. Aggarwal. *Outlier Analysis*. Springer New York, 2013.
- [2] Charu C. Aggarwal and Philip S. Yu. Outlier detection for high dimensional data. *SIGMOD Rec.*, 30:37–46, May 2001.
- [3] Andreas Arning, Rakesh Agrawal, and Prabhakar Raghavan. A linear method for deviation detection in large databases. In *KDD*, pages 164–169, 1996.
- [4] Stephen D. Bay and Mark Schwabacher. Mining distance-based outliers in near linear time with randomization and a simple pruning rule. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '03, pages 29–38, New York, NY, USA, 2003. ACM.
- [5] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. Lof: identifying density-based local outliers. In *ACM sigmod record*, volume 29, pages 93–104. ACM, 2000.
- [6] Glenn W. Brier. Verification of Forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1):1–3, January 1950.
- [7] Mujdat Cetin and William Clem Karl. Feature-enhanced synthetic aperture radar image formation based on nonquadratic regularization. *Image Processing, IEEE Transactions on*, 10(4):623–631, 2001.
- [8] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM Comput. Surv.*, 41(3):15:1–15:58, July 2009.
- [9] Krzysztof Dembczynski, Weiwei Cheng, and Eyke Hüllermeier. Bayes optimal multilabel classification via probabilistic classifier chains. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 279–286. Omnipress, 2010.
- [10] Pedro Garcia-Teodoro, J Diaz-Verdejo, Gabriel Maciá-Fernández, and Enrique Vázquez. Anomaly-based network intrusion detection: Techniques, systems and challenges. *computers & security*, 28(1):18–28, 2009.
- [11] Milos Hauskrecht, Iyad Batal, Michal Valko, Shyam Visweswaran, Gregory F Cooper, and Gilles Clermont. Outlier detection for patient monitoring and alerting. *Journal of Biomedical Informatics*, 46(1):47–55, 2013.
- [12] Milos Hauskrecht, Michal Valko, Branislav Kveton, Shyam Visweswaram, and Gregory Cooper. Evidence-based anomaly detection. In *Annual American Medical Informatics Association Symposium*, pages 319–324, November 2007.
- [13] D.M. Hawkins. *Identification of Outliers*. Monographs on applied probability and statistics. Chapman and Hall, 1980.
- [14] Victoria Hodge and Jim Austin. A survey of outlier detection methodologies. *Artif. Intell. Rev.*, 22(2):85–126, 2004.
- [15] Hans-Peter Kriegel, Peer Kröger, and Arthur Zimek. Outlier detection techniques. In *Tutorial at the 2010 SIAM International Conference on Data Mining*, 2010.
- [16] Abhishek Kumar, Shankar Vembu, Aditya Krishna Menon, and Charles Elkan. Learning and inference in probabilistic classifier chains with beam search. In *Proceedings of the 2012 European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer-Verlag, 2012.
- [17] Hancong Liu, Sirish Shah, and Wei Jiang. On-line outlier detection and data cleaning. *Computers & Chemical Engineering*, 28(9):1635–1647, 2004.
- [18] Andrew Y Ng. Feature selection,  $l_1$  vs.  $l_2$  regularization, and rotational invariance. In *Proceedings of the twenty-first international conference on Machine learning*. ACM, 2004.
- [19] Spiros Papadimitriou, Hiroyuki Kitagawa, Phillip B Gibbons, and Christos Faloutsos. Loci: Fast outlier detection using the local correlation integral. In *Data Engineering, 2003. Proceedings. 19th International Conference on*, pages 315–326. IEEE, 2003.
- [20] Jesse Read, Bernhard Pfahringer, Geoff Holmes, and Eibe Frank. Classifier chains for multi-label classification. In *Proceedings of the European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer-Verlag, 2009.
- [21] P. J. Rousseeuw and A. M. Leroy. *Robust Regression and Outlier Detection*. John Wiley & Sons, Inc., New York, NY, USA, 1987.
- [22] Peter J. Rousseeuw. Least median of squares regression. *Journal of the American Statistical Association*, 79(388):pp. 871–880, 1984.
- [23] Peter J. Rousseeuw and Bert C. van Zomeren. Unmasking multivariate outliers and leverage points. *Journal of the American Statistical Association*, 85(411):pp. 633–639, 1990.
- [24] Bernhard Schölkopf, Robert C Williamson, Alex J Smola, John Shawe-Taylor, and John C Platt. Support vector method for novelty detection. *NIPS*, 12:582–588, 1999.
- [25] Xiuyao Song, Mingxi Wu, Christopher Jermaine, and Sanjay Ranka. Conditional anomaly detection. *IEEE Trans. on Knowl. and Data Eng.*, 19(5):631–645, 2007.
- [26] Grigorios Tsoumakas, Ioannis Katakis, and Ioannis Vlahavas. Mining multi-label data. In *Data Mining and Knowledge Discovery Handbook*, pages 667–685. Springer US, 2010.
- [27] Shiguo Wang. A comprehensive survey of data mining-based accounting-fraud detection research. In *Intelligent Computation Technology and Automation (ICTA), 2010 International Conference on*, volume 1, pages 50–53, May 2010.
- [28] Cort J Willmott and Kenji Matsuura. Advantages of the mean absolute error (mae) over the root mean square error (rmse) in assessing average model performance. *Climate research*, 30(1):79–82, 2005.
- [29] Weng-Keen Wong, Andrew Moore, Gregory Cooper, and Michael Wagner. Bayesian network anomaly pattern detection for disease outbreaks. In *Proceedings of the Twentieth International Conference on Machine Learning*, pages 808–815. AAAI Press, August 2003.
- [30] M. Zhang and Z. Zhou. A review on multi-label learning algorithms. *Knowledge and Data Engineering, IEEE Transactions on*, PP(99):1, 2013.

- [31] Min-Ling Zhang and Kun Zhang. Multi-label learning by exploiting label dependency. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '10, pages 999–1008. ACM, 2010.